

CS 245: Database System Principles

Notes 03: Disk Organization

Hector Garcia-Molina

CS 245

Notes 3

1

Topics for today

- How to lay out data on disk
- How to move it to memory

CS 245

Notes 3

2

What are the data items we want to store?

- a salary
- a name
- a date
- a picture

⇒ What we have available: Bytes



CS 245

Notes 3

3

To represent:

- Integer (short): 2 bytes
e.g., 35 is

00000000 00100011

- Real, floating point
 n bits for mantissa, m for exponent...

CS 245

Notes 3

4

To represent:

- Characters

→ various coding schemes suggested,
most popular is ascii

Example:

A: 1000001
a: 1100001
5: 0110101
LF: 0001010

CS 245

Notes 3

5

To represent:

- Boolean

e.g., TRUE 1111 1111
FALSE 0000 0000

- Application specific
e.g., RED → 1 GREEN → 3
BLUE → 2 YELLOW → 4 ...

⇒ Can we use less than 1 byte/code?
Yes, but only if desperate...

CS 245

Notes 3

6

To represent:

- Dates
 - e.g.: - Integer, # days since Jan 1, 1900
 - 8 characters, YYYYMMDD
 - 7 characters, YYYYDDD
(not YYMMDD! Why?)
- Time
 - e.g. - Integer, seconds since midnight
 - characters, HHMMSSFF

CS 245

Notes 3

7

To represent:

- String of characters

- Null terminated

e.g.,

c	a	t	⊗		
---	---	---	---	--	--

- Length given

e.g.,

3	c	a	t	⊗	
---	---	---	---	---	--

- Fixed length

CS 245

Notes 3

8

To represent:

- Bag of bits

Length	Bits
--------	------

CS 245

Notes 3

9

Key Point

- Fixed length items
- Variable length items
 - usually length given at beginning

CS 245

Notes 3

10

Also

- Type of an item: Tells us how to interpret
(plus size if fixed)

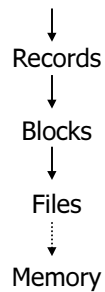
CS 245

Notes 3

11

Overview

Data Items



CS 245

Notes 3

12

Record - Collection of related data items (called FIELDS)

E.g.: Employee record:
name field,
salary field,
date-of-hire field, ...

CS 245

Notes 3

13

Types of records:

- Main choices:
 - FIXED vs VARIABLE FORMAT
 - FIXED vs VARIABLE LENGTH

CS 245

Notes 3

14

Fixed format

A SCHEMA (not record) contains following information

- # fields
- type of each field
- order in record
- meaning of each field

CS 245

Notes 3

15

Example: fixed format and length

Employee record

- (1) E#, 2 byte integer
- (2) E.name, 10 char.
- (3) Dept, 2 byte code

} Schema

55	s	m	i	t	h							02
83	j	o	n	e	s							01

} Records

CS 245

Notes 3

16

Variable format

- Record itself contains format "Self Describing"

CS 245

Notes 3

17

Example: variable format and length

2	5	I	46	4	S	4	F	O	R	D
---	---	---	----	---	---	---	---	---	---	---

Fields ↑
Code identifying field as E# ↑
Integer type ↑
Code for Ename ↑
String type ↑
Length of str. ↑

Field name codes could also be strings, i.e. TAGS

CS 245

Notes 3

18

Variable format useful for:

- "sparse" records
- repeating fields
- evolving formats

.....→ But may waste space...

CS 245

Notes 3

19

- EXAMPLE: var format record with repeating fields
Employee → one or more → children

3	E_name: Fred	Child: Sally	Child: Tom
---	--------------	--------------	------------

CS 245

Notes 3

20

Note: Repeating fields does not imply
- variable format, nor
- variable size

John	Sailing	Chess	--
------	---------	-------	----

- Key is to allocate maximum number of repeating fields (if not used → null)

CS 245

Notes 3

21

☆ Many variants between fixed - variable format:

Ex. #1: Include record type in record

5	27	...
---	----	-----

↑ record type
tells me what to expect
(i.e. points to schema)

← record length

CS 245

Notes 3

22

Record header - data at beginning that describes record

May contain:

- record type
- record length
- time stamp
- other stuff ...

CS 245

Notes 3

23

Ex #2 of variant between FIXED/VAR format

- Hybrid format
- one part is fixed, other variable

E.g.: All employees have E#, name, dept other fields vary.

25	Smith	Toy	2	Hobby:chess	retired
----	-------	-----	---	-------------	---------

↑
of var fields

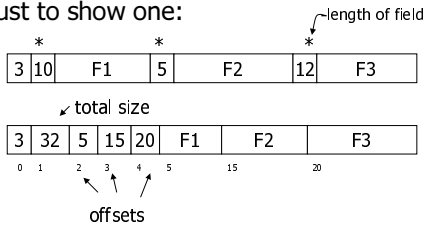
CS 245

Notes 3

24

☆ Also, many variations in internal organization of record

Just to show one:



CS 245

Notes 3

25

Question:

We have seen examples for

- * Fixed format and length records
- * Variable format and length records

(a) Does fixed format and variable length make sense?

(b) Does variable format and fixed length make sense?

CS 245

Notes 3

26

Other interesting issues:

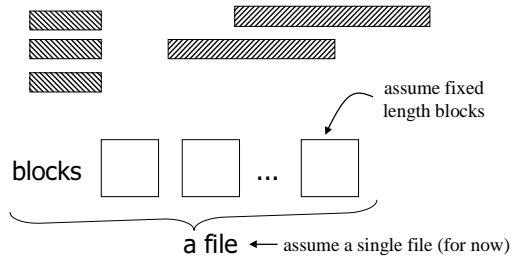
- Compression
 - within record - e.g. code selection
 - collection of records - e.g. find common patterns
- Encryption

CS 245

Notes 3

27

Next: placing records into blocks



CS 245

Notes 3

28

Options for storing records in blocks:

- (1) separating records
- (2) spanned vs. unspanned
- (3) mixed record types – clustering
- (4) split records
- (5) sequencing
- (6) indirection

CS 245

Notes 3

29

(1) Separating records



(a) no need to separate - fixed size recs.

(b) special marker

(c) give record lengths (or offsets)

- within each record
- in block header

CS 245

Notes 3

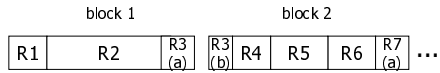
30

(2) Spanned vs. Unspanned

- Unspanned: records must be within one block



- Spanned

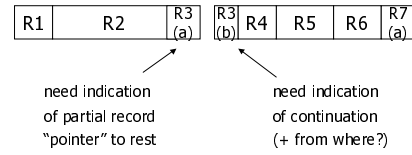


CS 245

Notes 3

31

With spanned records:



need indication
of partial record
"pointer" to rest

need indication
of continuation
(+ from where?)

CS 245

Notes 3

32

Spanned vs. unspanned:

- Unspanned is much simpler, but may waste space...
- Spanned essential if
record size > block size

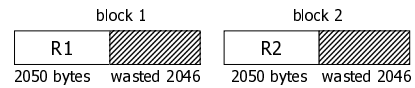
CS 245

Notes 3

33

Example

10^6 records
each of size 2,050 bytes (fixed)
block size = 4096 bytes



- Total wasted = 2×10^9 Utiliz = 50%
- Total space = 4×10^9

CS 245

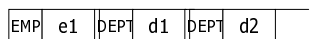
Notes 3

34

(3) Mixed record types

- Mixed - records of different types
(e.g. EMPLOYEE, DEPT)
allowed in same block

e.g., a block



CS 245

Notes 3

35

Why do we want to mix?

Answer: CLUSTERING

Records that are frequently
accessed together should be
in the same block

CS 245

Notes 3

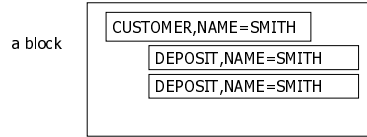
36

Compromise:

No mixing, but keep related records in same cylinder ...

Example

Q1: select A#, C_NAME, C_CITY, ...
from DEPOSIT, CUSTOMER
where DEPOSIT.C_NAME =
CUSTOMER.C.NAME



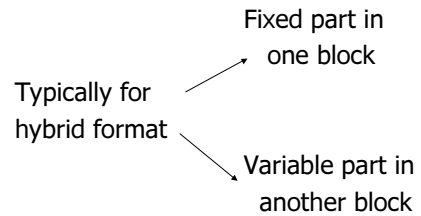
- If Q1 frequent, clustering good
- But if Q2 frequent

Q2: SELECT *

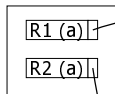
FROM CUSTOMER

CLUSTERING IS COUNTER PRODUCTIVE

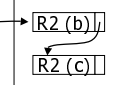
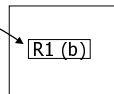
(4) Split records



Block with fixed recs.



Block with variable recs.



This block also has fixed recs.

Question

What is difference between
- Split records
- Simply using two different record types?

(5) Sequencing

- Ordering records in file (and block) by some key value

Sequential file (\Rightarrow sequenced)

CS 245

Notes 3

43

Why sequencing?

Typically to make it possible to efficiently read records in order

(e.g., to do a merge-join — discussed later)

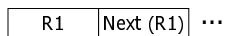
CS 245

Notes 3

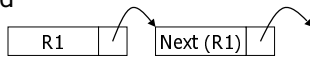
44

Sequencing Options

(a) Next record physically contiguous



(b) Linked



CS 245

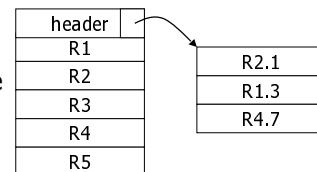
Notes 3

45

Sequencing Options

(c) Overflow area

Records
in sequence



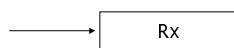
CS 245

Notes 3

46

(6) Indirection

- How does one refer to records?



Many options:

Physical \longleftrightarrow Indirect

CS 245

Notes 3

47

☆ Purely Physical

E.g., Record Address or ID = $\left. \begin{array}{l} \text{Device ID} \\ \text{Cylinder \#} \\ \text{Track \#} \\ \text{Block \#} \\ \text{Offset in block} \end{array} \right\} \text{Block ID}$

CS 245

Notes 3

48

☆ Fully Indirect

E.g., Record ID is arbitrary bit string

CS 245 Notes 3 49

Tradeoff

Flexibility \longleftrightarrow Cost
 to move records of indirection
 (for deletions, insertions)

CS 245 Notes 3 50

Physical \longleftrightarrow Indirect

↑
 Many options
 in between ...

CS 245 Notes 3 51

Ex #1 Indirection in block

CS 245 Notes 3 52

Block header - data at beginning that describes block

May contain:

- File ID (or RELATION or DB ID)
- This block ID
- Record directory
- Pointer to free space
- Type of block (e.g. contains recs type 4; is overflow, ...)
- Pointer to other blocks "like it"
- Timestamp ...

CS 245 Notes 3 53

Ex. #2 Use logical block #'s understood by file system

REC ID \rightarrow File ID
 Block #
 Record # or Offset

CS 245 Notes 3 54

File system map may be "Semi-physical"...

File F1: physical address of block 1
table of bad blocks:

{	B57	→	XXX
	B107	→	YYY

Rest can be computed via formula...

CS 245 Notes 3 55

Num. Blocks: 20
Start Block: 1000
Block Size: 100
Bad Blocks:
 3 → 20,000
 7 → 15,000

Where is Block # 2?
Where is Block # 3?

File DEFINITION

CS 245 Notes 3 56

Options for storing records in blocks

- (1) Separating records
- (2) Spanned vs. Unspanned
- (3) Mixed record types - Clustering
- (4) Split records
- (5) Sequencing
- (6) Indirection

CS 245 Notes 3 57

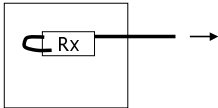
Other Topics

- (1) Insertion/Deletion
- (2) Buffer Management
- (3) Comparison of Schemes

CS 245 Notes 3 58

Deletion

Block



The diagram shows a large rectangle labeled 'Block'. Inside it, on the left side, is a smaller rectangle labeled 'Rx'. An arrow points from the right side of the 'Rx' rectangle to the right edge of the 'Block' rectangle.

CS 245 Notes 3 59

Options:

- (a) Immediately reclaim space
- (b) Mark deleted
 - May need chain of deleted records (for re-use)
 - Need a way to mark:
 - special characters
 - delete field
 - in map

CS 245 Notes 3 60

☆ As usual, many tradeoffs...

- How expensive is to move valid record to free space for immediate reclaim?
- How much space is wasted?
 - e.g., deleted records, delete fields, free space chains,...

CS 245

Notes 3

61

Concern with deletions

Dangling pointers



CS 245

Notes 3

62

Solution #1: Do not worry

CS 245

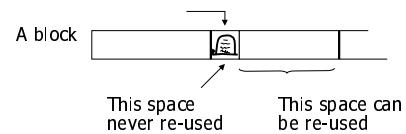
Notes 3

63

Solution #2: Tombstones

E.g., Leave "MARK" in map or old location

- Physical IDs



CS 245

Notes 3

64

Solution #2: Tombstones

E.g., Leave "MARK" in map or old location

- Logical IDs

map	
ID	LOC
7788	

Never reuse ID 7788 nor space in map...

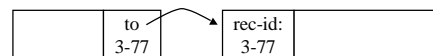
CS 245

Notes 3

65

Solution #3 (?):

- Place record ID within every record
- When you follow a pointer, check if it leads to correct record



Does this work???

If space reused, won't new record have same ID?

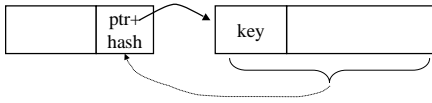
CS 245

Notes 3

66

Solution #4 (?):

- To point, use (pointer + hash)
or (pointer + key)?



- What if record modified???

CS 245

Notes 3

67

Insert

Easy case: records not in sequence

- Insert new record at end of file or in deleted slot
- If records are variable size, not as easy...

CS 245

Notes 3

68

Insert

Hard case: records in sequence

- If free space "close by", not too bad...
- Or use overflow idea...

CS 245

Notes 3

69

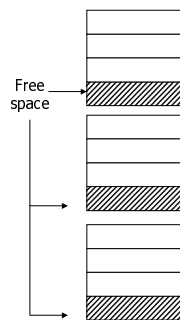
Interesting problems:

- How much free space to leave in each block, track, cylinder?
- How often do I reorganize file + overflow?

CS 245

Notes 3

70



CS 245

Notes 3

71

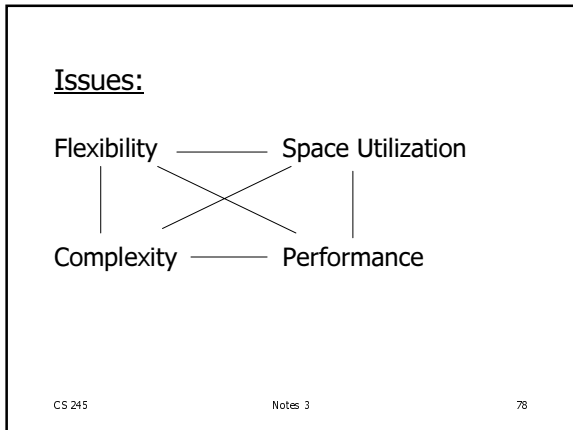
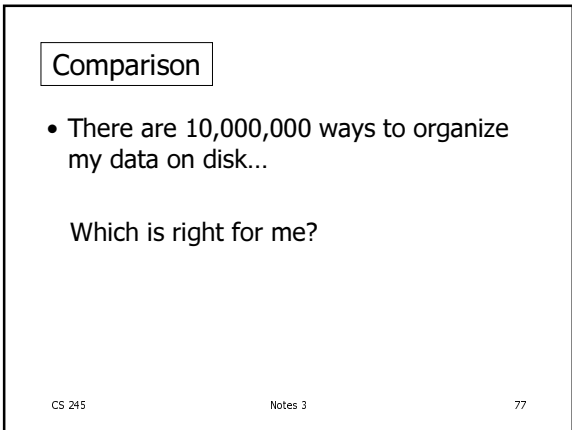
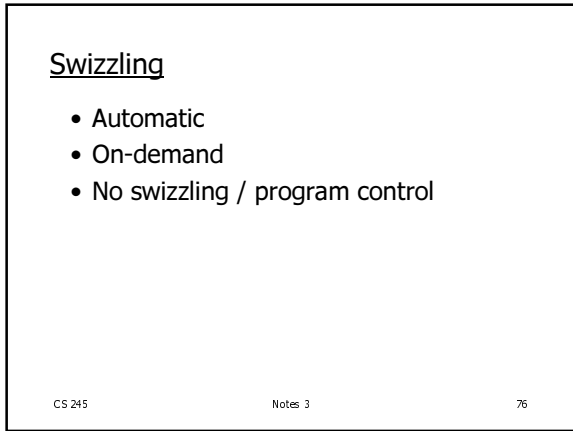
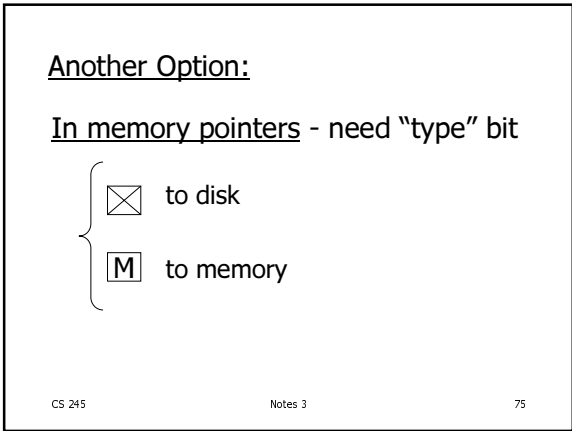
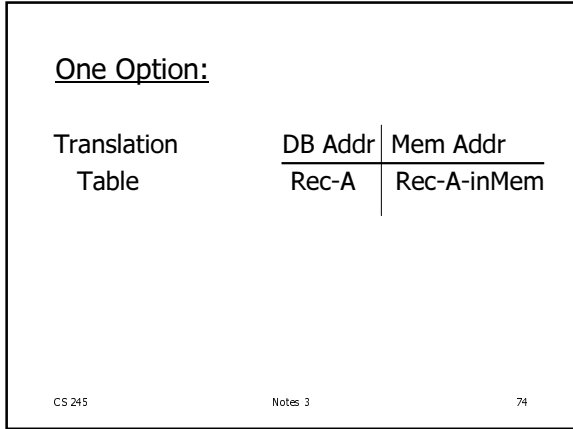
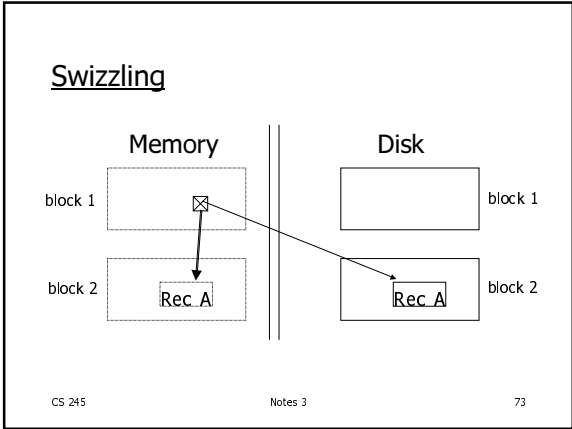
Buffer Management

- DB features needed
 - Pinned blocks
 - Forced output (flush)
 - Why LRU may be bad
 - Double buffering
 - Swizzling
- } Read Textbook!
..... in Notes02

CS 245

Notes 3

72



☆ To evaluate a given strategy, compute following parameters:

-> space used for expected data

-> expected time to

- fetch record given key
- fetch record with next key
- insert record
- append record
- delete record
- update record
- read all file
- reorganize file

CS 245

Notes 3

79

Example

How would you design Megatron 3000 storage system? (for a relational DB, low end)

- Variable length records?
- Spanned?
- What data types?
- Fixed format?
- Record IDs ?
- Sequencing?
- How to handle deletions?

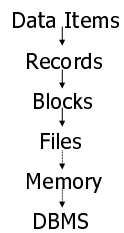
CS 245

Notes 3

80

Summary

- How to lay out data on disk



CS 245

Notes 3

81

Next

How to find a record quickly,
given a key

CS 245

Notes 3

82